# New Dimensions in Microarchitecture

**Harnessing 3D Integration Technologies**

**Kerry Bernstein**
**IBM T.J. Watson Research Center**
**Yorktown Heights, NY**

**6 March, 2007**     **San Jose, California**
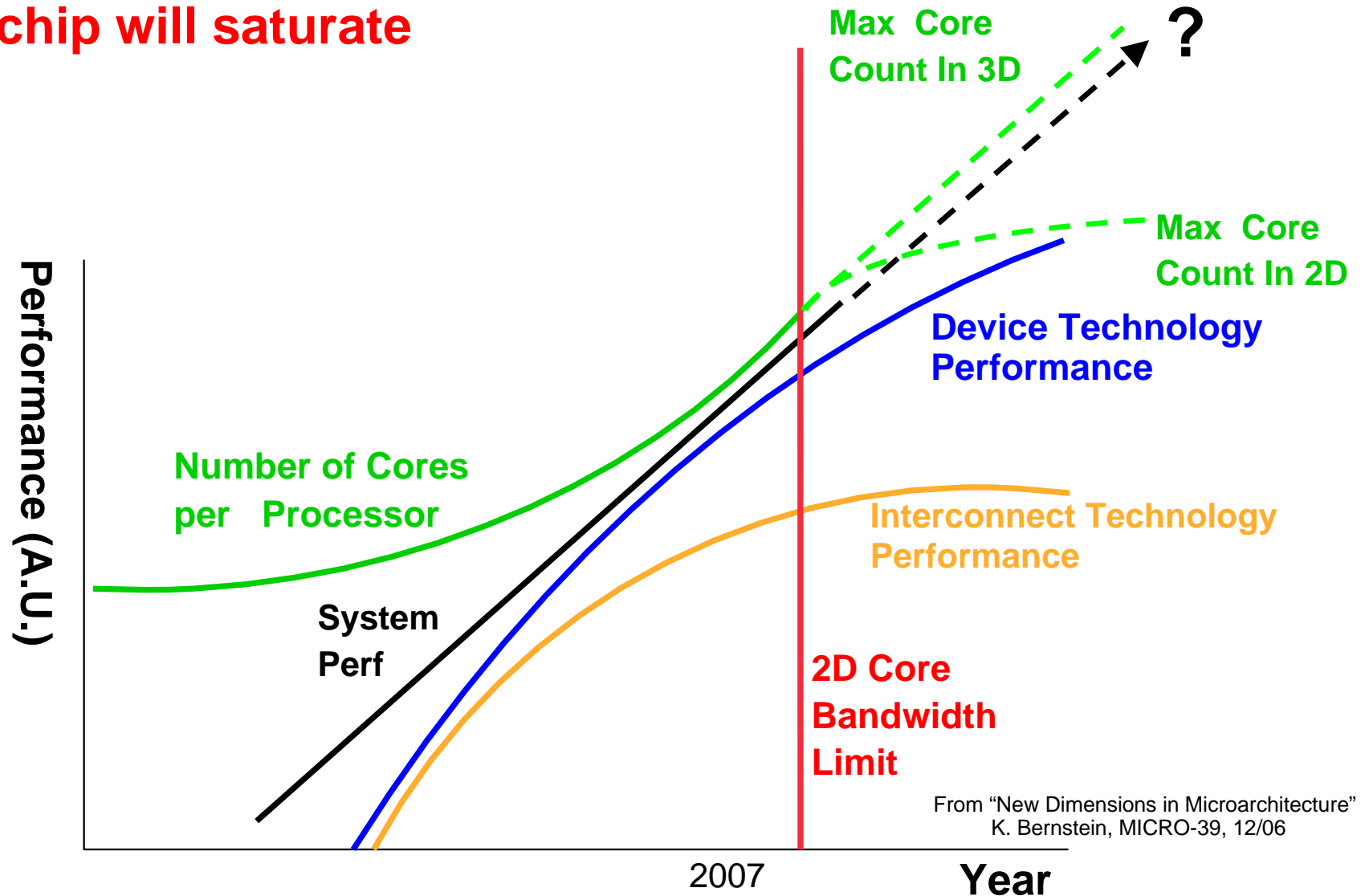
DARPA Microsystems Technology Symposium                "Escher Envy" courtesy of David Bryant

| 1. REPORT DATE | 2. REPORT TYPE | 3. DATES COVERED |
|---|---|---|
| **06 MAR 2007** | **N/A** | **-** |

| 4. TITLE AND SUBTITLE | 5a. CONTRACT NUMBER |
|---|---|
| **New Dimensions in Microarchitecture Harnessing 3D Integration Technologies** | 5b. GRANT NUMBER |
| | 5c. PROGRAM ELEMENT NUMBER |

| 6. AUTHOR(S) | 5d. PROJECT NUMBER |
|---|---|
| | 5e. TASK NUMBER |
| | 5f. WORK UNIT NUMBER |

| 7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) | 8. PERFORMING ORGANIZATION REPORT NUMBER |
|---|---|
| **IBM Corporation** | |

| 9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) | 10. SPONSOR/MONITOR'S ACRONYM(S) |
|---|---|
| | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) |

12. DISTRIBUTION/AVAILABILITY STATEMENT

**Approved for public release, distribution unlimited**

13. SUPPLEMENTARY NOTES

**DARPA Microsystems Technology Symposium held in San Jose, California on March 5-7, 2007. Presentations, The original document contains color images.**

14. ABSTRACT

15. SUBJECT TERMS

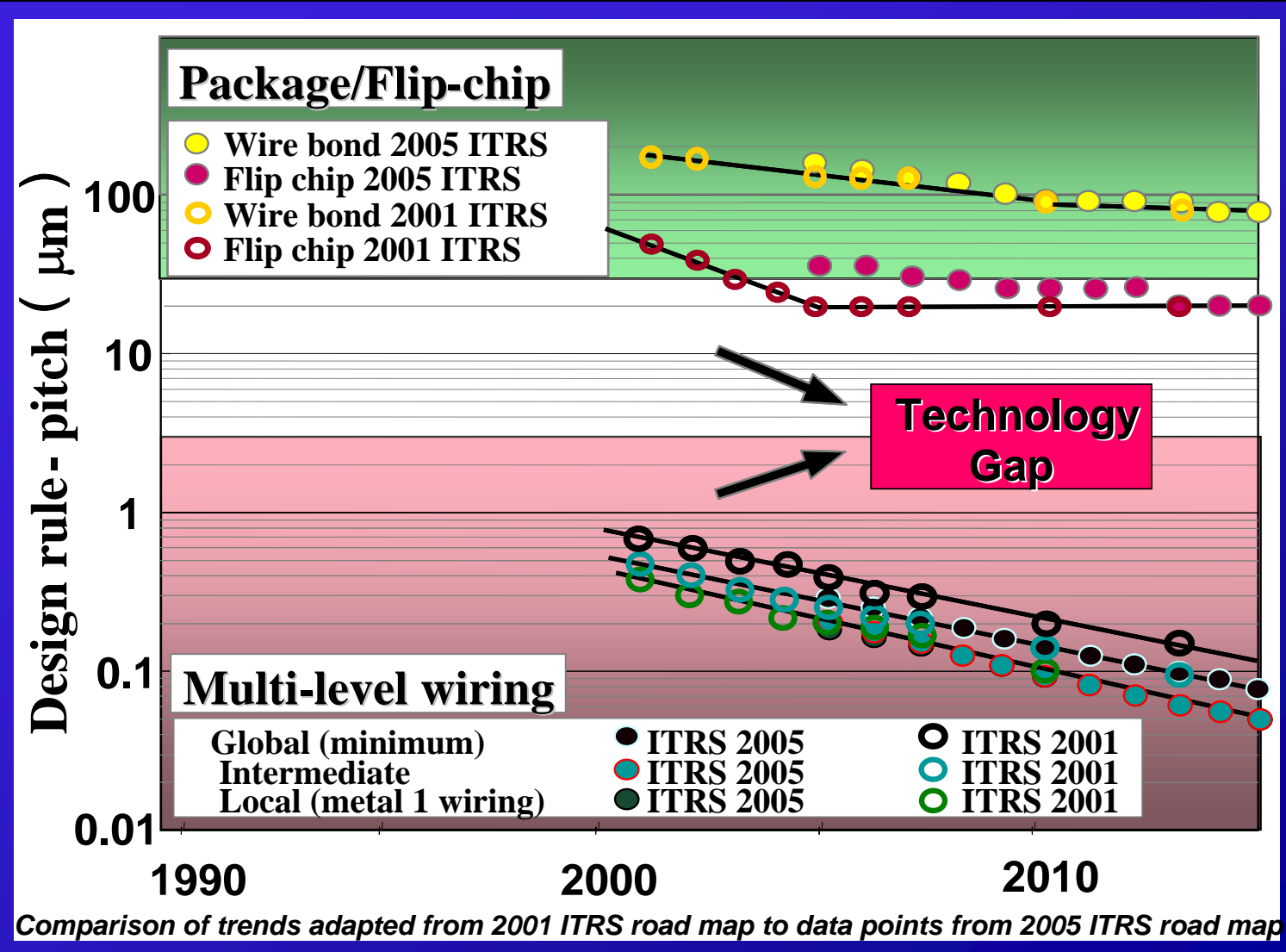| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES | 19a. NAME OF RESPONSIBLE PERSON |
|---|---|---|---|---|---|
| a. REPORT **unclassified** | b. ABSTRACT **unclassified** | c. THIS PAGE **unclassified** | **UU** | **12** | |

# Server Trends

- Frequency no longer increasing

    - Logic speed scaled faster than memory bus

    - (Processor clocks / Bus clock) consumes bandwidth

- More speculation; attempts to prefetch

    - Wrong guesses increase miss traffic

- Shortening linesize limited by directory as cache size grows

    - But doubling linesize doubles bus occupancy

- Cores / die increasing each generation

    - Multiplies off-chip bus transactions by N / 2*Sqrt(2)

- More threads per core, and increase in virtualization

    - Multiplies off-chip bus transactions by N

- Processors / SMP increasing

    - Aggravates queueing throughout the system

**Without more bandwidth at low latencies, core counts on chip will saturate**

Max Core Count In 3D

?

Max Core Count In 2D

Device Technology Performance

Performance (A.U.)

Number of Cores per Processor

Interconnect Technology Performance

System Perf

2D Core Bandwidth Limit

From "New Dimensions in Microarchitecture"
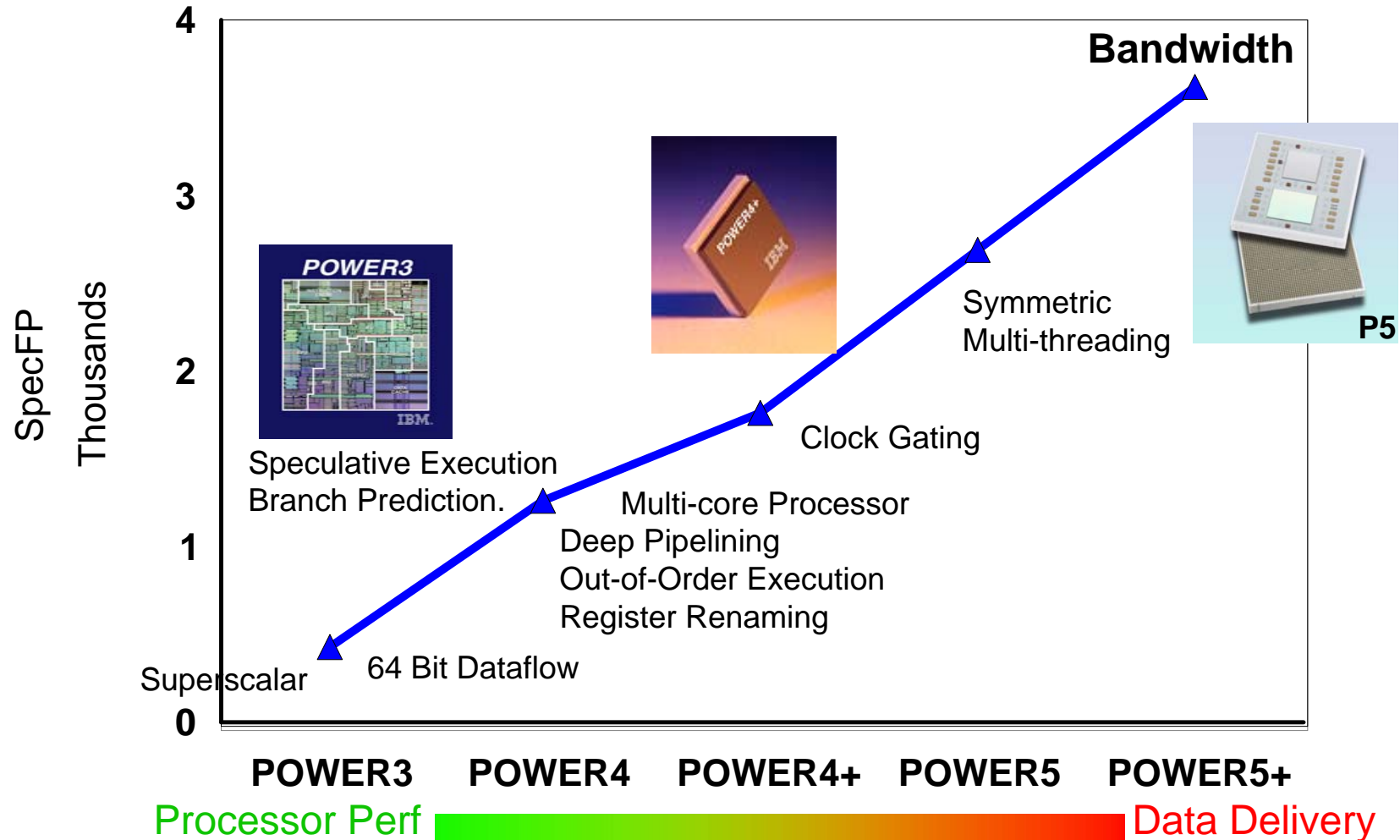K. Bernstein, MICRO-39, 12/06

2007

Year

**3D extends transfer of performance from the device to the core level**

# Chip-Package Technology Gap

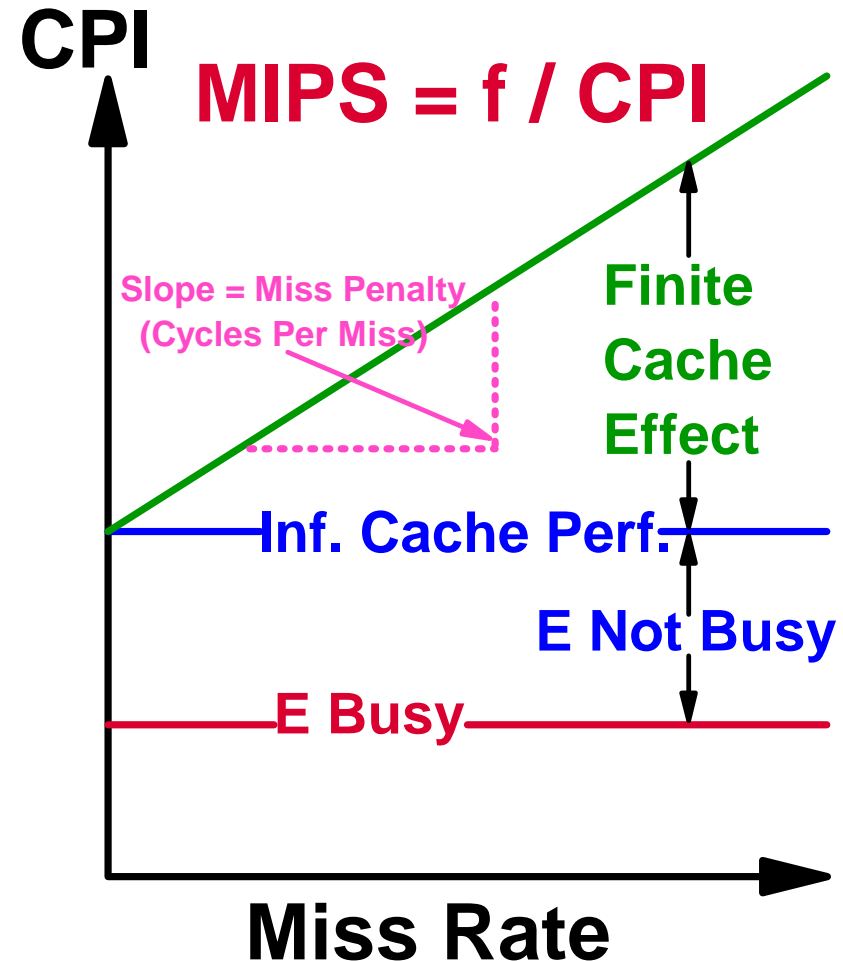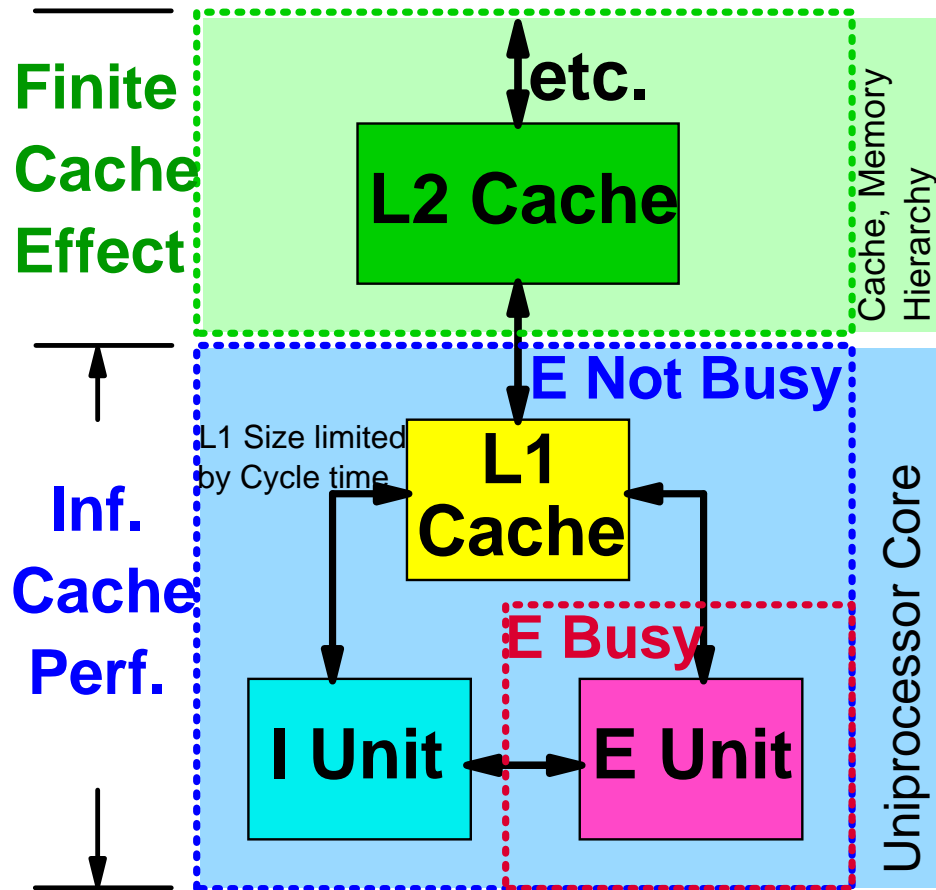**Technology gap in the design rule between on-chip wiring and packaging interconnects**

DARPA MTS · March 6, 2007 · © 2006 IBM Corporation

# POWER Series Architectural Perf Contributions



POWER3, POWER4, POWER4+, POWER5, POWER5+ performance chart showing SpecFP (Thousands) on the vertical axis (0 to 4).

- **Bandwidth**
- Symmetric Multi-threading
- Clock Gating
- Multi-core Processor
- Deep Pipelining
- Out-of-Order Execution
- Register Renaming
- Speculative Execution Branch Prediction.
- 64 Bit Dataflow
- Superscalar

**P5**

Processor Perf ▸ **Transaction Rate Dependence** ◂ Data Delivery

New Dimensions in Microarchitecture

# Components of Processor Performance

**Finite Cache Effect**

**Inf. Cache Perf.**

Cache, Memory Hierarchy

etc.

**L2 Cache**

E Not Busy

L1 Size limited by Cycle time

**L1 Cache**

E Busy

**I Unit**

**E Unit**

Uniprocessor Core

**CPI**

**MIPS = f / CPI**

Slope = Miss Penalty (Cycles Per Miss)

**Finite Cache Effect**

**Inf. Cache Perf.**

**E Not Busy**

**E Busy**

**Miss Rate**

Delay is sequentially determined by a) ideal processor, b) access to local cache, and c) refill of cache

6

From ISCA '06
Keynote address by
Phil Emma, IBM

# Queueing Effects vs. Log Miss Rate

From ISCA '06 Keynote address by Phil Emma, IBM

# What Is Bandwidth Used For?

**In a computer, it is mostly for handling cache misses:**

Miss

Access

Processor Events

Bus Events

Leading Edge

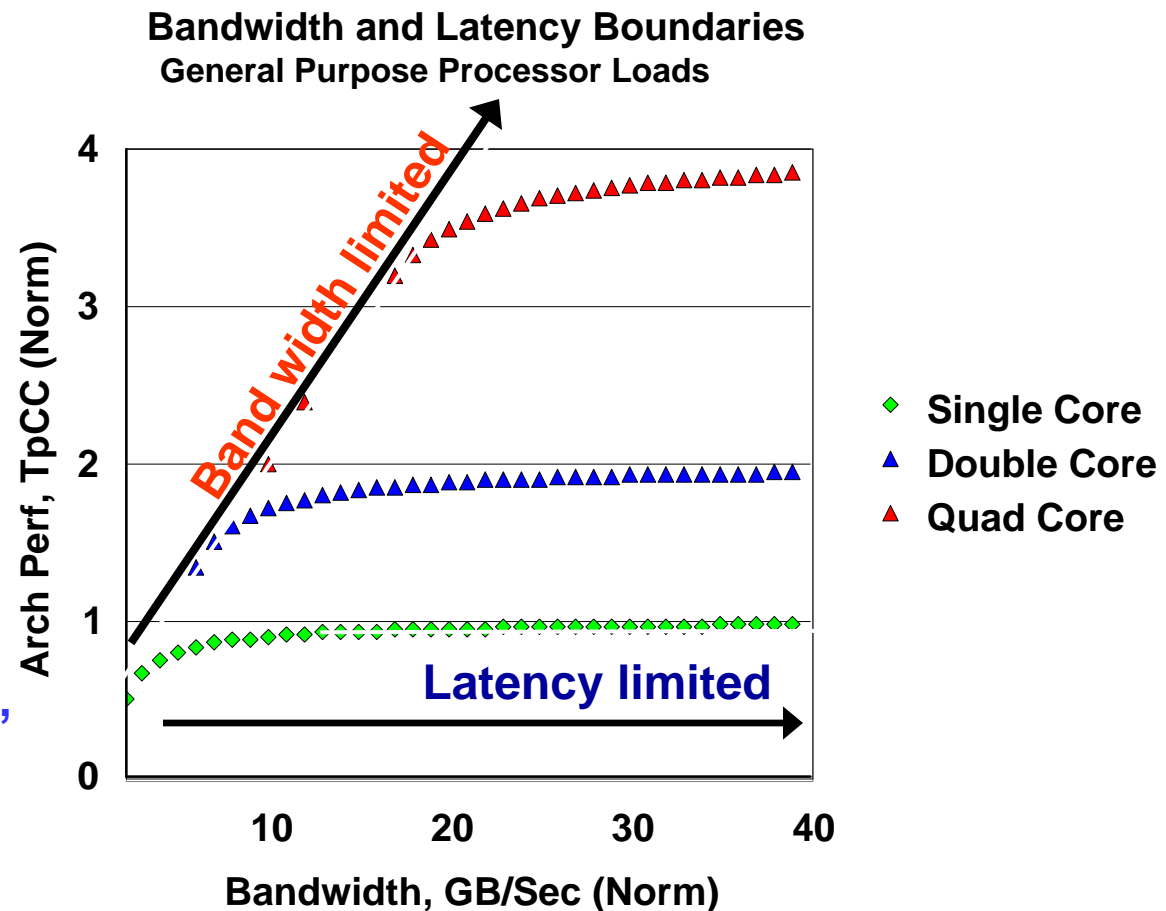Rest of Cache Line

Trailing Edge

First Data

Last Data

Time

**Miss Penalty = Leading Edge + Effects(Trailing Edge)**

# 3D - Bandwidth and Latency

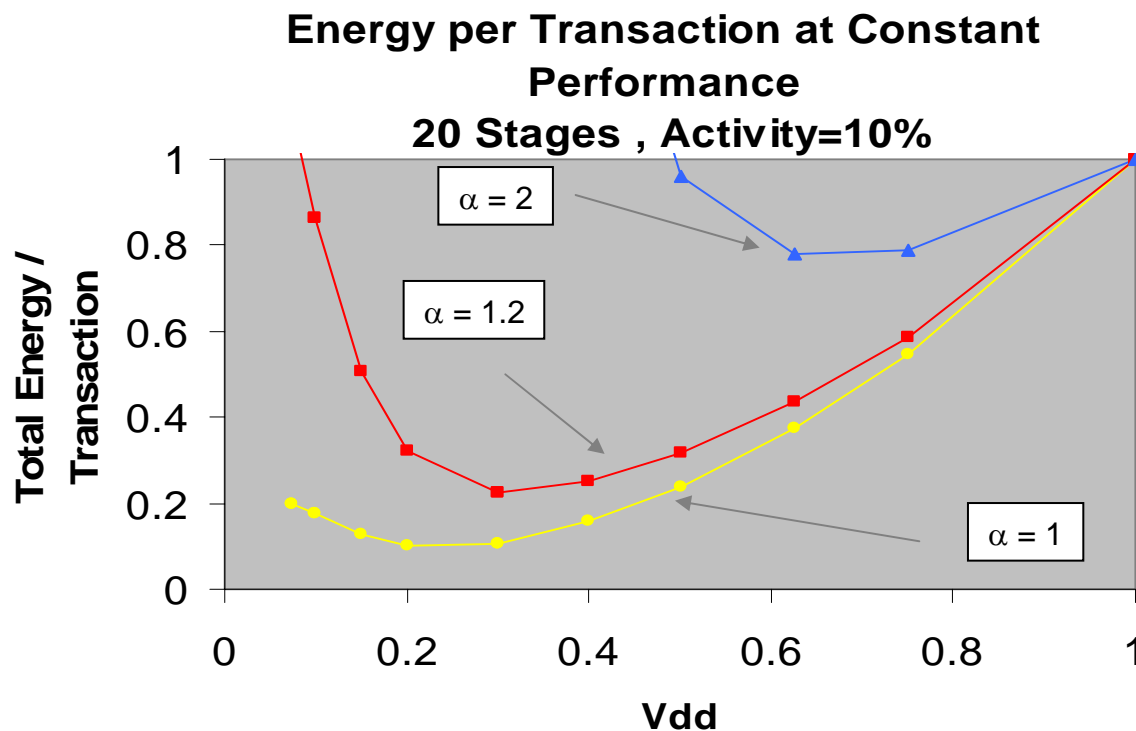**Processor load trade-off between I/O Bandwidth, Bus Latency.**

**- For generic workloads, uni-processor perf saturates bandwidth benefit, becomes latency-limited.**

**- As core counts increase, I/O Bandwidth becomes increasingly important**

**Bandwidth and Latency Boundaries**
**General Purpose Processor Loads**



Band width limited

Latency limited

◆ **Single Core**
▲ **Double Core**
▲ **Quad Core**

**Arch Perf, TpCC (Norm)**

**Bandwidth, GB/Sec (Norm)**

**3D opportunity for improving High Perf Compute thru-put in sustaining a higher number of cores per chip**

# Low Vdd Technology and Parallelism

- Energy optimum for fixed performance as function of $V_{dd}$, $V_T$ and effectiveness of parallelism

  - $\alpha$ determines the device (circuit) over head to maintain constant performance through parallelism

  - $\alpha = 1$ **no overhead**: half the speed double the devices

  - $\alpha > 1$ **increasing overhead**: passive power becomes dominant

**Energy per Transaction at Constant Performance
20 Stages , Activity=10%**



$$P \sim P_0 \frac{d_0}{d}\left(\frac{N}{N_0}\right)^{\frac{1}{\alpha}}$$

N=number of ckts,
d=ckt delay
$N_o$= number of ckts at 1V
$d_o$= ckt delay at 1V
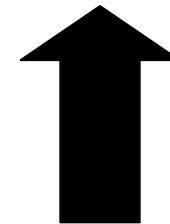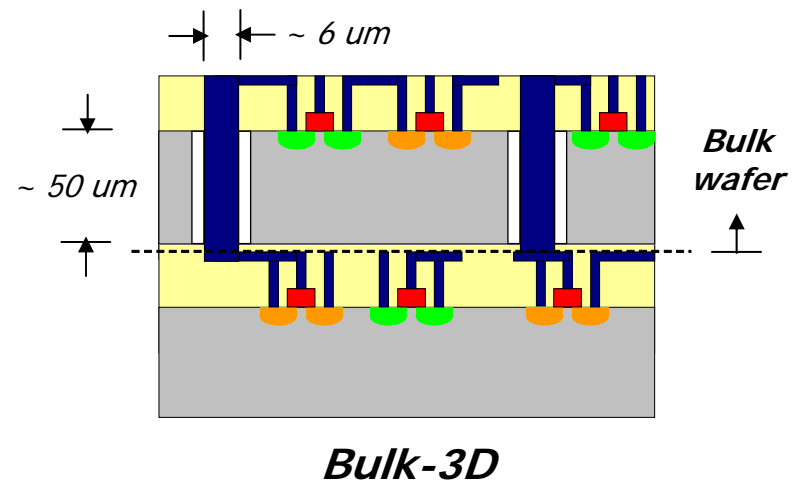
# Two Classes of 3DI Processes at IBM



SOI-3D

Bulk-3D

**SOI top layer**
*Advantage: Smallest 3D vias*

**Bulk top layer**
*Advantage: Broader foundry compatibility*

# Summary

- $\lambda$P architecture tricks to avoid atomistic, QM scaling boundaries overwhelm present interconnect technologies

- Integration into Z-plane again postpones interconnect-related limitations to extending classic scaling.

- No aspect of architecture or technology remains 2D, so why even view chips as being monolithic anymore?

- Transaction retirement rate dependence on data delivery is *increasing*: dependence on $\lambda$P performance and CMOS device speed is *decreasing*

- 3D Integration improves storage density and access to that storage

- The last remaining opportunity in CMOS to save power is in delivery of data rather than in its generation.